

#### Plan

#### Today

- Wrap-up neural nets (for now)
- Regularization
  - (Make sure they aren't too powerful <sup>(2)</sup>)
  - Regularization with L2 norm
  - Regularization optimization
  - Regularization with L1 norm

# Regularization with L2 norm

Example: Linear regression with polynomial features



Which is model do you prefer, assuming both have zero training error?

Model structure (for both models):

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \theta_5 x^5 + \theta_6 x^6 + \theta_7 x^7 + \theta_8 x^8$$

Model parameters:

$$\boldsymbol{\theta} = [\theta_0, \ \theta_1, \ \theta_2, \ \theta_3, \ \theta_4, \ \theta_5, \ \theta_6, \ \theta_7, \ \theta_8]^T$$

**A.** 
$$\theta_A = [-190.0, -135.0, 310.0, 45.0, -62.0, 90.0, -82.0, -40.0, 29.0]^T$$

**B.** 
$$\theta_B = [25.5, -6.4, -0.8, 0.0, 6.6, -4.4, 0.2, -2.9, 0.1]^T$$

#### Which is model do you prefer, assuming both have zero training error?

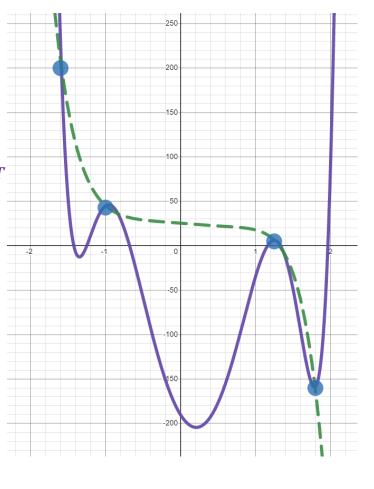
Model structure (for both models):

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \theta_5 x^5 + \theta_6 x^6 + \theta_7 x^7 + \theta_8 x^8$$

Model parameters:

$$\boldsymbol{\theta} = [\theta_0, \ \theta_1, \ \theta_2, \ \theta_3, \ \theta_4, \ \theta_5, \ \theta_6, \ \theta_7, \ \theta_8]^T$$

$$\boldsymbol{\theta}_A = [-190.0, -135.0, 310.0, 45.0, -62.0, 90.0, -82.0, -40.0, 29.0]^T$$
 $\boldsymbol{\theta}_B = [25.5, -6.4, -0.8, 0.0, 6.6, -4.4, 0.2, -2.9, 0.1]^T$ 



# Overfitting

Definition: The problem of **overfitting** is when the model captures the noise in the training data instead of the underlying structure

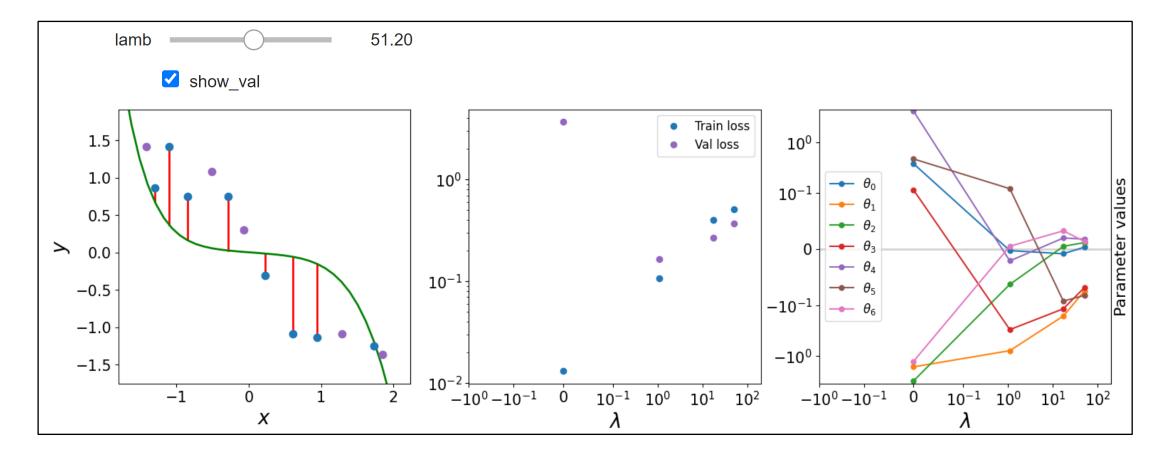
#### Overfitting can occur in all the models we've seen so far:

- Decision Trees (e.g. when tree is too deep)
- K-NN (e.g. when k is small)
- Linear Regression (e.g. with nonlinear features or extraneous features)
- Logistic Regression (e.g. with nonlinear features or extraneous features)
- Neural networks

#### Best of both worlds

How can we keep the expressive power of a complex model while still avoiding overfitting?

Notebook demo: regression regularization.ipynb

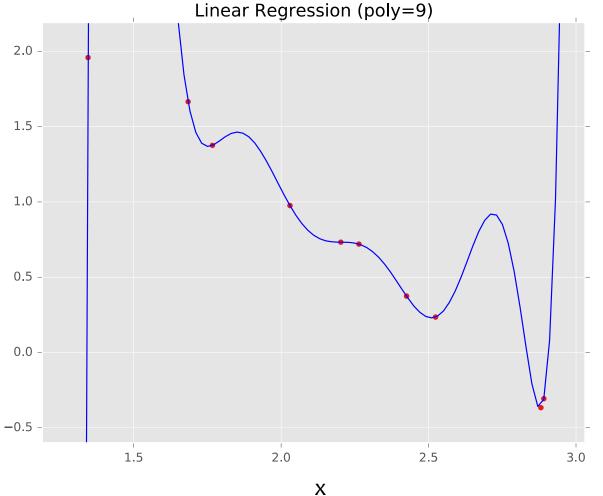


# Example: Linear Regression

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + \mathbf{b}$  where f(.) is a polynomial basis function

У	х	x <sup>2</sup>	 x <sup>9</sup>	
2.0	1.2	(1.2)2	 (1.2)9	
1.3	1.7	(1.7)2	 (1.7)9	
0.1	2.7	(2.7)2	 (2.7)9	y
1.1	1.9	(1.9)2	 (1.9)9	•

true "unknown"
target function is
linear with
negative slope
and gaussian
noise



## Symptoms of Overfitting

	M=0	M = 1	M = 3	M = 9
$\overline{\theta_0}$	0.19	0.82	0.31	0.35
$ heta_1$		-1.27	7.99	232.37
$ heta_2$			-25.43	-5321.83
$ heta_3$			17.37	48568.31
$ heta_4$				-231639.30
$ heta_5$				640042.26
$ heta_6$				-1061800.52
$ heta_7$				1042400.18
$ heta_8$				-557682.99
$ heta_9$				125201.43

Motivation: Regularization

Occam's Razor: prefer the simplest hypothesis

#### What does it mean for a hypothesis (or model) to be simple?

- 1. small number of features (model selection)
- 2. small number of "important" features (feature reduction)
- $\rightarrow$  3. small values for associated parameters  $1/\beta/(\frac{1}{2})$  22-norw

#### Key idea:

Define regularizer  $r(\theta)$  that we will add to our minimization objective to keep the model simple.

#### $r(\theta)$ should be:

- Small for a simple model
- Large for a complex model

L2 norm: square-root of sum of squares

$$r(\theta) = ||\theta||_2^2$$

L1 norm: sum of absolute values

LO norm: count of non-zero values

$$f(\theta) = \|\boldsymbol{\theta}\|_2^2$$

**A.** 
$$\theta_A = [6, 3, -4, -2]^T$$

**B.** 
$$\theta_B = [0, 3, -4, 0]^T$$

Which model do you prefer?

$$\boldsymbol{\theta}_A = [-190.0, -135.0, 310.0, 45.0]^T$$
 Training error: 0.0

 $\boldsymbol{\theta}_B = [0.0, 0.0, 0.0, 0.0]^T$  Training error: 34.2

Notebook demo: regression regularization.ipynb on course website

What is the best value for lambda?  $\theta_0 + \theta_1 x_1 + \theta_1 x_1^2$  $\hat{\boldsymbol{\theta}} = \operatorname{argmin} \boldsymbol{J}(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$ Train loss 1.5 10<sup>0</sup> 1.0  $10^{-1}$ 0.5  $10^{-1}$ 0.0 -0.5 $-10^{-1}$ -1.0 $-10^{0}$ -1.5 $10^{1}$   $10^{2}$  $-10^{0} - 10^{-1}$  $-10^{0} - 10^{-1}$  $10^{-1}$   $10^{0}$  $10^{-1}$ Х

Notebook demo: regression regularization.ipynb on course website

What is the best value for lambda?  $\hat{\boldsymbol{\theta}} = \operatorname{argmin} \boldsymbol{J}(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$ Train loss 1.5  $10^{0}$ Val loss 1.0  $10^{0}$  $10^{-1}$ 0.5 0.0  $10^{-1}$ -0.5 $-10^{-1}$ -1.0 $-10^{0}$ -1.5 $10^{-2}$  $10^1 10^2$  $-10^{0} - 10^{-1}$  $-10^{0} - 10^{-1}$  $10^{-1}$   $10^{0}$  $10^{1}$ Х

Given objective function:  $J(\theta) = \frac{|\vec{y} - \vec{x}\vec{\theta}|/2}{2}$ 

Goal is to find: 
$$\hat{\boldsymbol{\theta}} = \operatorname*{argmin} J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

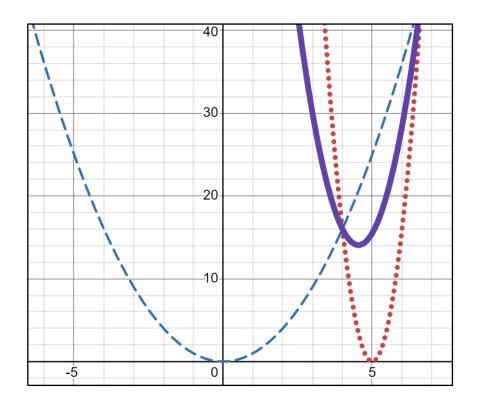
**Key idea**: Define regularizer  $r(\theta)$  s.t. we tradeoff between fitting the data and keeping the model simple

Choose form of 
$$r(\theta)$$
:  $\neq ||\theta||_2^7$ 

#### L2 Demos

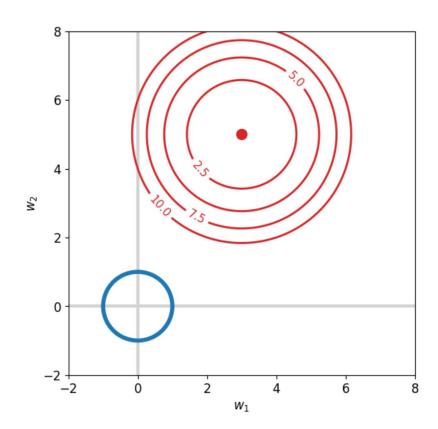
Desmos: 1-D

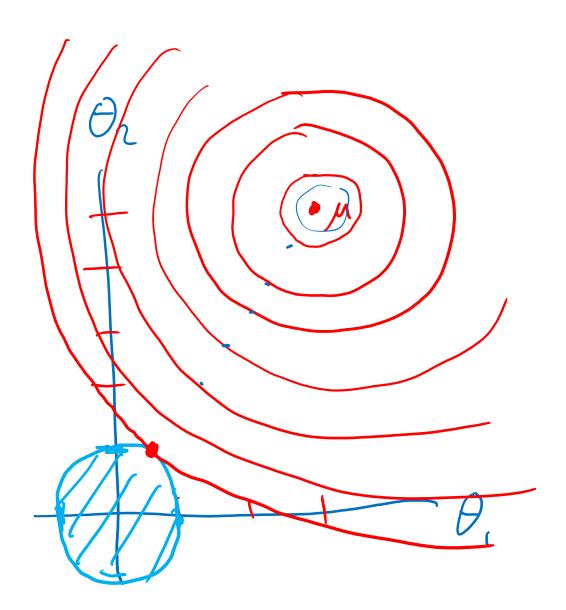
**Regularization Interpolation** 



Notebook: 2-D

L1 sparsity.ipynb (L2 part for now)





$$J(\theta_{i},\theta_{i}) = ||\vec{\theta} - \vec{n}|| \qquad M = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$
min 
$$J(\theta_{i},\theta_{i})$$

$$\theta$$
s.t. 
$$||\theta||_{2} \leq |$$

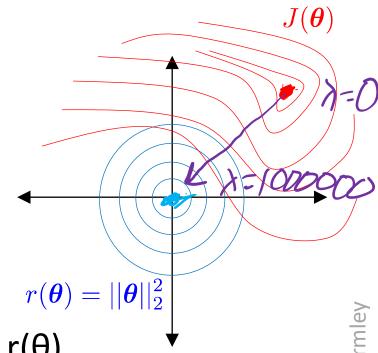
Suppose we are minimizing  $J'(\theta)$  where

$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

As  $\lambda$  increases, the minimum of J'( $\theta$ ) will...



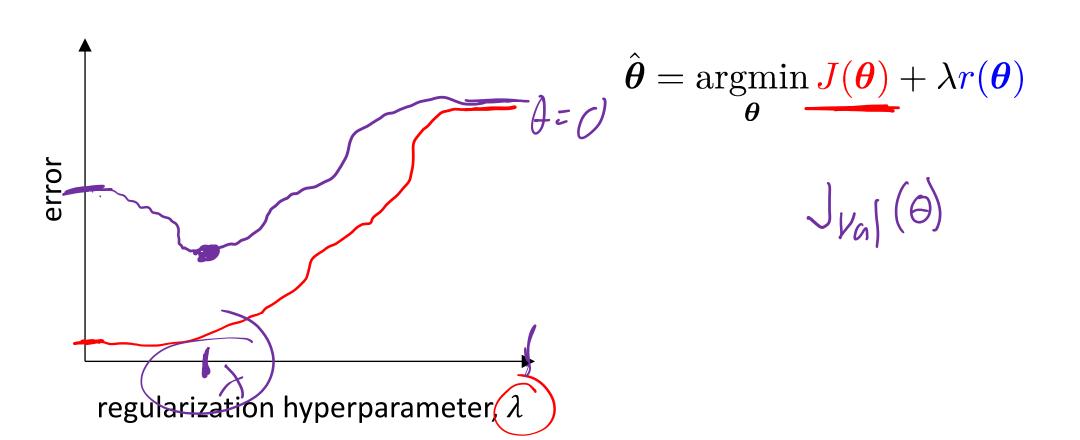
- B. ...move towards the minimum of  $J(\theta)$
- (C.)...move towards the minimum of  $r(\theta)$ 
  - D. ...move towards a theta vector of positive infinities
  - E. ...move towards a theta vector of negative infinities
  - F. ...stay the same



## Regularization Exercise

#### In-class Exercise

- 1. Plot train error vs. regularization hyperparameter (cartoon)
- 2. Plot validation error vs. regularization hyperparameter (cartoon)

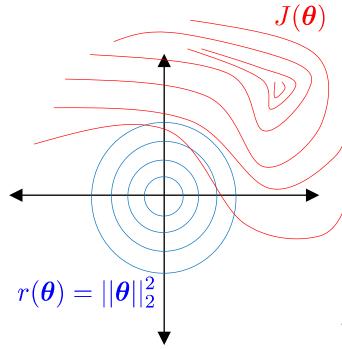


Suppose we are minimizing  $J'(\theta)$  where

$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

As we increase  $\lambda$  from zero, the **validation** error will...  $r(\theta) = ||\theta||_2^2$ 

- A. ...increase
- B. ...decrease
- C. ...first increase, then decrease
- (D.)...first decrease, then increase
  - E. ...stay the same



As we increase  $\lambda$ , our model is more likely to:

- A. Overfit
- B. Underfit

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$



 $X = (X_{raw} - M)/\alpha$ 

#### Don't Regularize the Bias (Intercept) Parameter

- In our models so far, the bias / intercept parameter is usually denoted by  $heta_0$  that is, the parameter for which we fixed  $x_0=1$
- Regularizers always avoid penalizing this bias / intercept parameter
- Why? Because otherwise the learning algorithms wouldn't be invariant to a shift in the y-values

#### Whitening Data

- It's common to whiten each feature by subtracting its mean and dividing by its variance
- For regularization, this helps all the features be penalized in the same units (e.g. convert both centimeters and kilometers to z-scores)

# Regularization Optimization

# Linear Regression with L2 Regularization

a.k.a Ridge regression or Tychonov regression

denom 
$$J(\theta) = ||y - X0||_{2}^{2} + \lambda ||\theta||_{2}^{2}$$
  
 $MxI \frac{\partial J}{\partial \theta} = 0 = -JX^{T}(\bar{y} - X\bar{\theta})^{-} + J\lambda \theta$   
 $= -X^{T}y + X^{T}X\theta + \lambda \theta$   
 $X^{T}y = (X^{T}X + \lambda I_{m})^{T}X^{T}y = \theta$ 

# Linear Algebra Timeout

$$A_V + cV = (A + cI)_V$$

Distribution of multiplication and addition with scalar involved

Origina | Broken | 
$$\begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 1 & 1 & 5 \\ 5 & 1 & 5 \end{bmatrix} + 3 \begin{bmatrix} 2 & 1 & 1 \\ 5 & 1 & 5 \\ 7 & 1 & 1 & 5 \end{bmatrix} + 2 \begin{bmatrix} 2 & 1 & 1 \\ 4 & 4 & 1 & 1 \\ 4 & 4 & 1 & 1 & 5 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 7 & 1 & 1 &$$

alar involved
$$Fixed$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 4 \end{bmatrix} + 3\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

$$= \begin{bmatrix} 13 \\ 22 \end{bmatrix}$$

# Regularization with L1 norm

#### Model Preference

Which is model do you prefer, assuming both have zero training error?

Model structure (for both models):

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \theta_5 x_5 + \theta_6 x_6 + \theta_7 x_7 + \theta_8 x_8$$

Model parameters:

$$\boldsymbol{\theta} = [\theta_0, \ \theta_1, \ \theta_2, \ \theta_3, \ \theta_4, \ \theta_5, \ \theta_6, \ \theta_7, \ \theta_8]^T$$

**A.** 
$$\theta_A = [-190.0, -135.0, 310.0, 45.0, -62.0, 90.0, -82.0, -40.0, 29.0]^T$$

$$\rightarrow B$$
,  $\theta_B = [25.5, -6.4, -0.8, 0.0, 6.6, -4.4, 0.2, -2.9, 0.1]^T$ 

What if **x** was a vector of input feature measurements (rather than polynomial features)?

## Motivation: Regularization

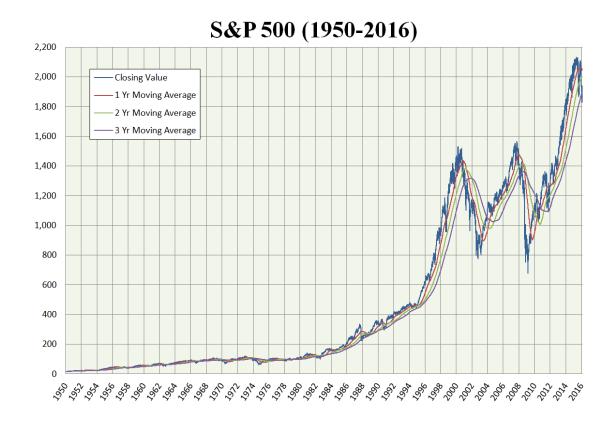
**Example: Stock Prices** 

Suppose we wish to predict Google's stock price at time t+1

What features should we use? (putting all computational concerns aside)

- Stock prices of all other stocks at times t, t-1, t-2, ..., t k
- Mentions of Google with positive / negative sentiment words in all newspapers and social media outlets

Do we believe that **all** of these features are going to be useful?



#### Key idea:

Define regularizer  $r(\theta)$  that we will add to our minimization objective to keep the model simple.

#### $r(\theta)$ should be:

- Small for a simple model
- Large for a complex model

$$=||\theta||_{1}=\sum_{j=1}^{N}|\theta_{j}|$$

$$-||\theta||_{0}=\sum_{j=1}^{N}|\theta_{j}|$$

$$-||\theta||_{0}=\sum_{j=1}^{N}|\theta_{j}|$$

$$\|\boldsymbol{\theta}\|_2^{\mathcal{E}}$$

$$\|oldsymbol{ heta}\|_1$$

$$\|\boldsymbol{\theta}\|_0$$

**A.** 
$$\theta_A = [6, 3, -4, -2]^T$$

**B.** 
$$\theta_B = [0, 3, -4, 0]^T$$

**Given** objective function:  $J(\theta)$ 

Goal is to find:  $\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta}) + \lambda \underline{r(\boldsymbol{\theta})}$ 

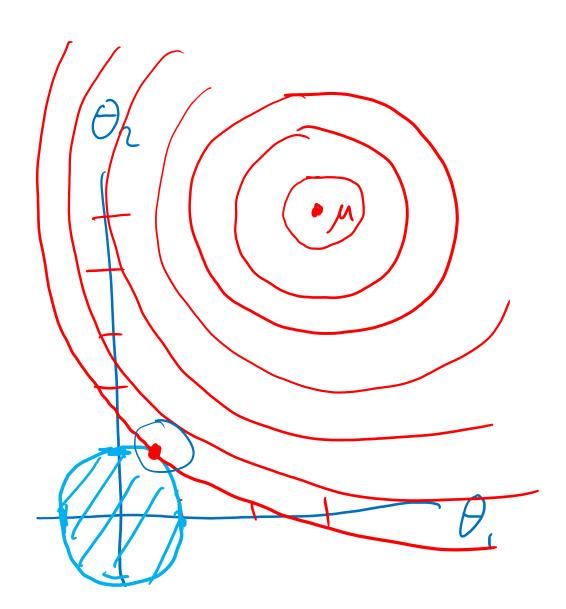
**Key idea**: Define regularizer  $r(\theta)$  s.t. we tradeoff between fitting the data and keeping the model simple

#### Choose form of $r(\theta)$ :

Example: q-norm (usually p-norm)

$$r(\boldsymbol{\theta}) = ||\boldsymbol{\theta}||_q = \left[\sum_{m=1}^{M} ||\theta_m||^q\right]^{(\frac{1}{q})}$$

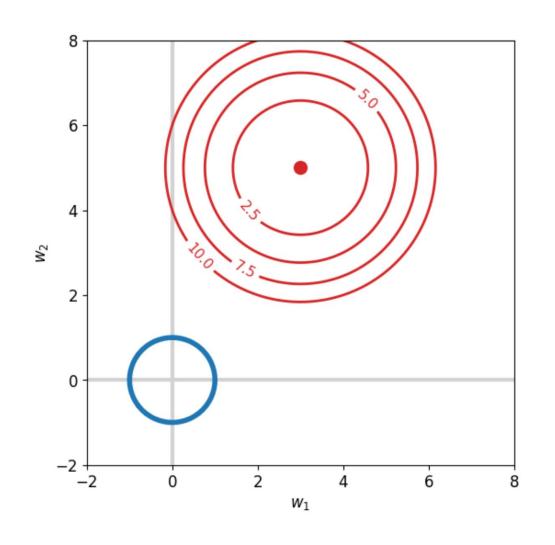
$\overline{q}$	$r(oldsymbol{ heta})$	yields parame- ters that are	name	optimization notes
0	$  \boldsymbol{\theta}  _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computa- tional solutions
$\begin{array}{c} 1 \\ 2 \end{array}$	$  oldsymbol{ heta}  _1 = \sum   heta_m  \ (  oldsymbol{ heta}  _2)^2 = \sum  heta_m^2$	zero values small values	L1 reg. L2 reg.	subdifferentiable differentiable

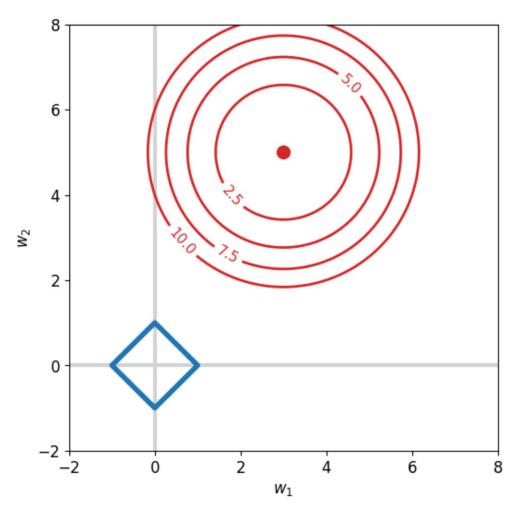


$$J(\theta, \theta_1) = ||\vec{\theta} - \vec{n}|| \qquad M = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$
min 
$$J(\theta, \theta_1)$$

$$\theta$$
s.t. 
$$||\theta||_2 \leq 1$$

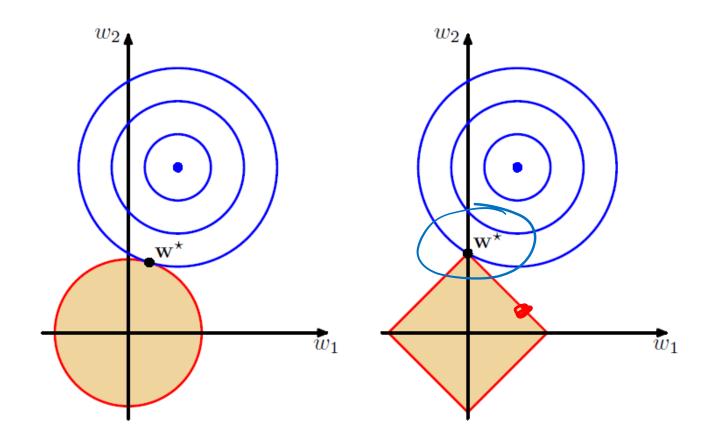
## L1 demo: L1 sparsity.ipynb





# L2 vs L1 Regularization

Combine original objective with penalty on parameters



Figures: Bishop, Ch 3.1.4

# L2 vs L1: Housing Price Example

Predict housing price from several features

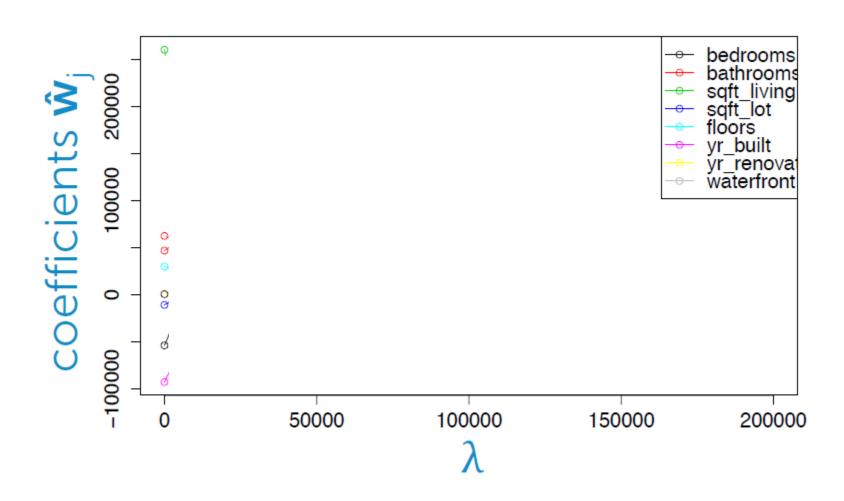
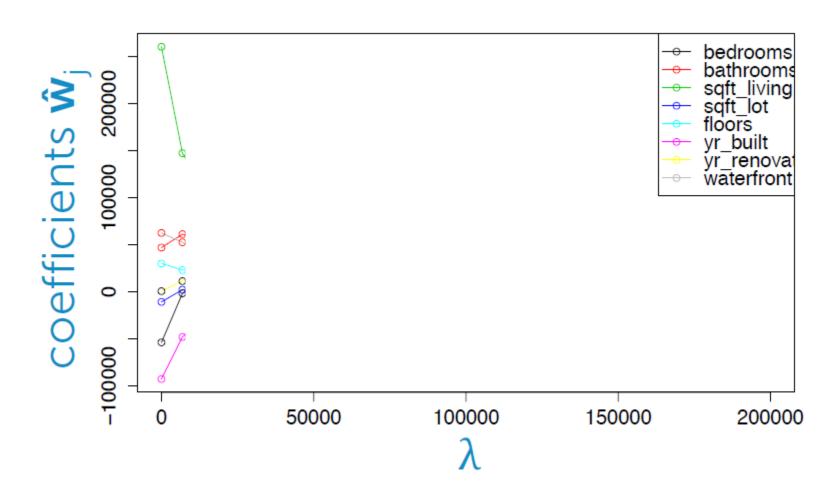
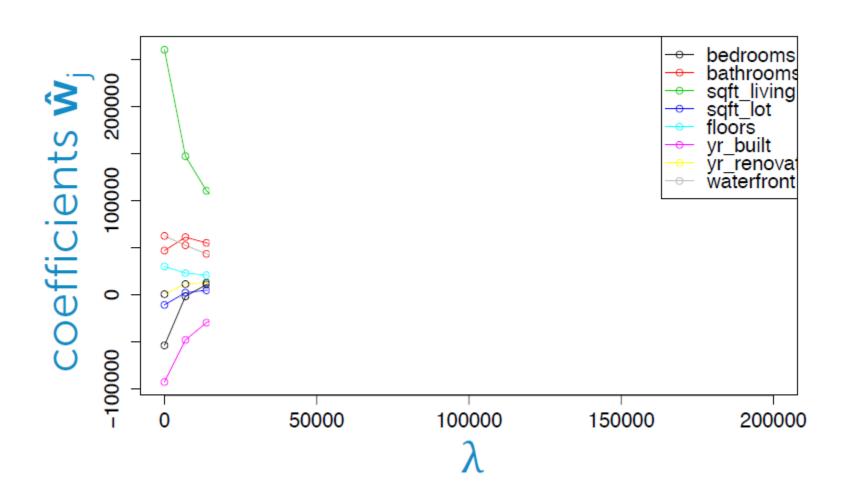


Figure: Emily Fox, University of Washington

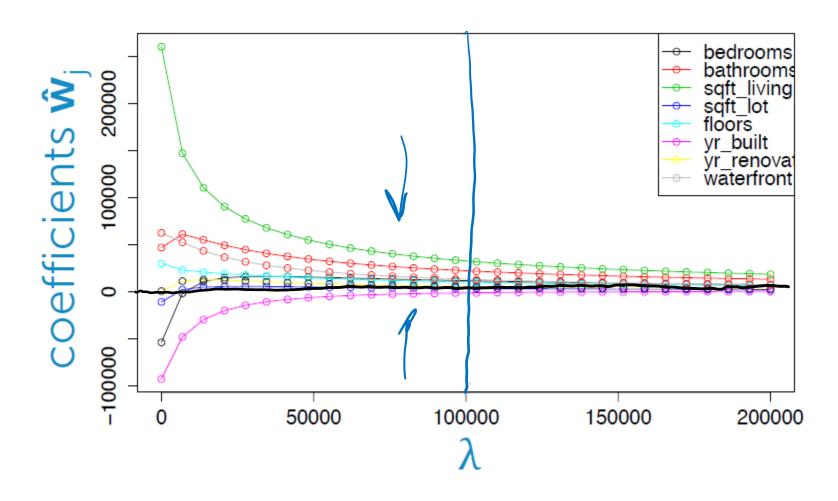
Predict housing price from several features



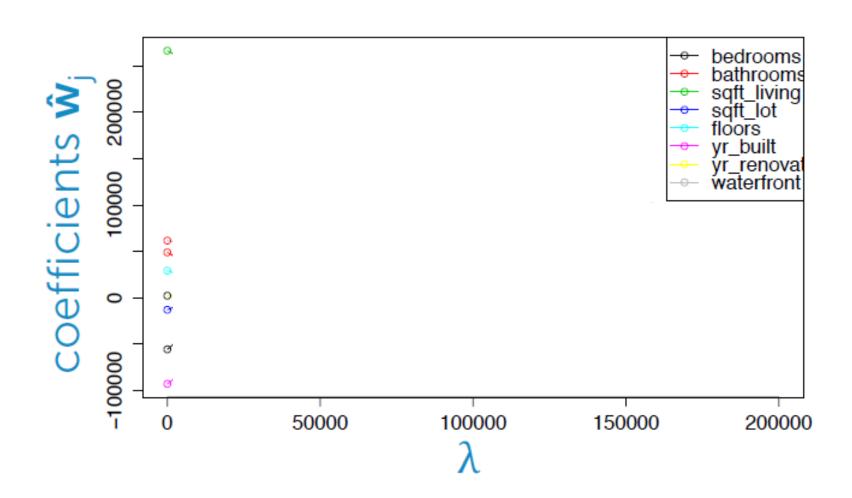
Predict housing price from several features



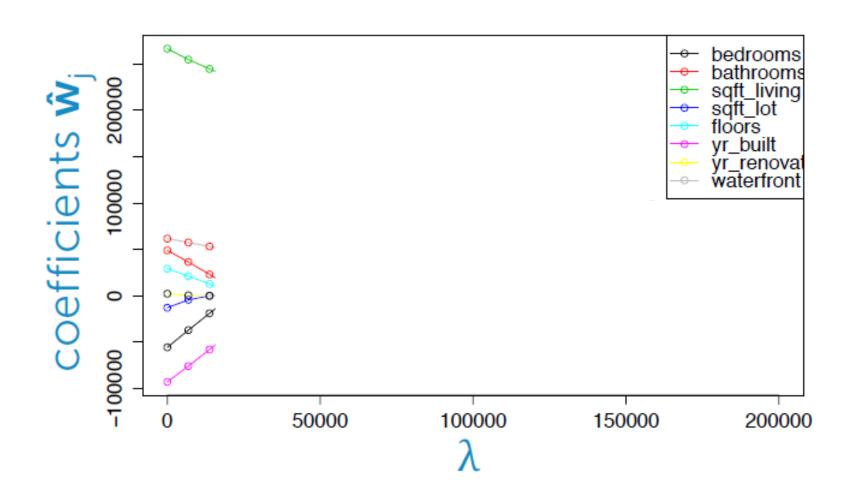
#### Predict housing price from several features



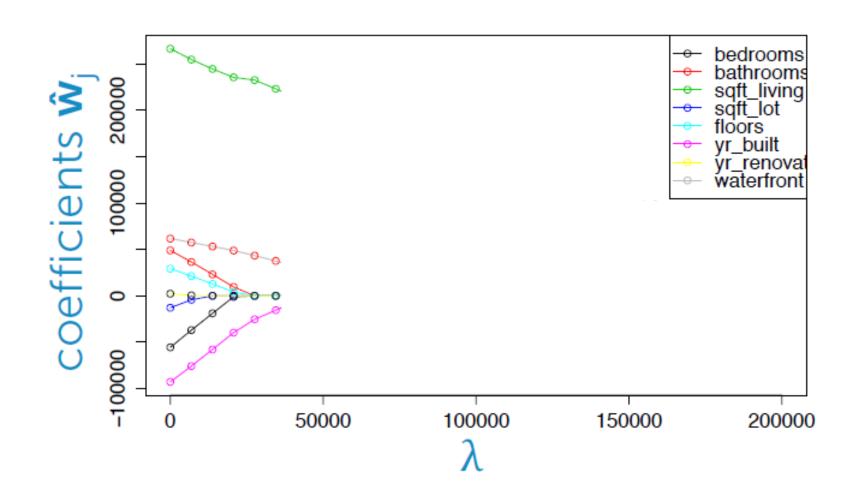
Predict housing price from several features



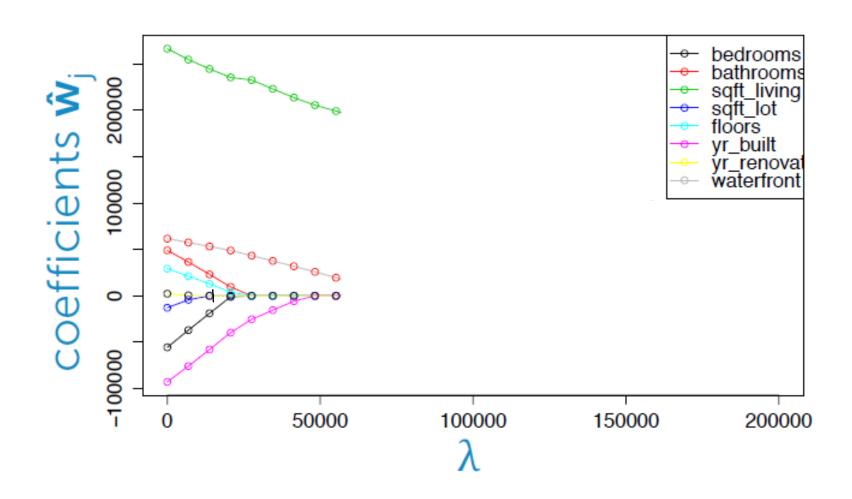
Predict housing price from several features



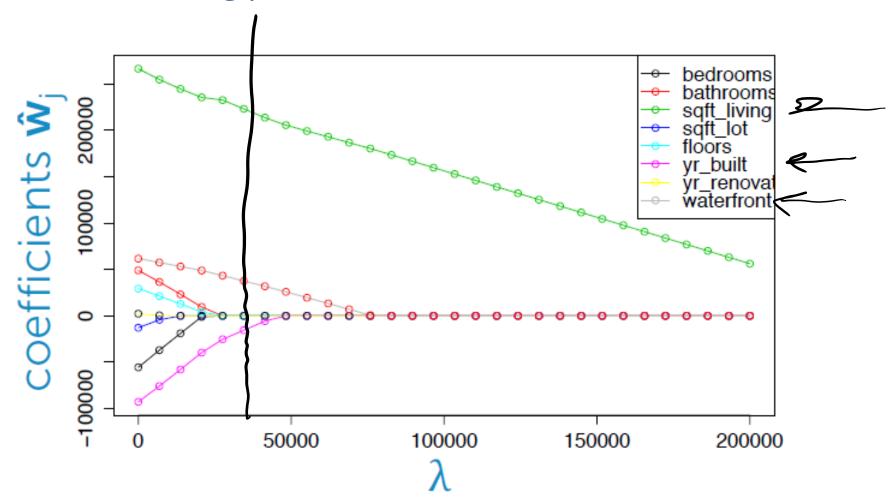
Predict housing price from several features



Predict housing price from several features



Predict housing price from several features



# Regularization as MAP

L1 and L2 regularization can be interpreted as maximum a-posteriori (MAP) estimation of the parameters

To be discussed later in the course...

#### Additional Slides

#### Logistic Regression with Nonlinear Features

Jupyter notebook demo



For this example, we construct **nonlinear features** (i.e. feature engineering)

Specifically, we add polynomials up to order 9 of the two original features  $x_1$  and  $x_2$ 

Thus our classifier is linear in the high-dimensional feature space, but the decision boundary is nonlinear when visualized in low-dimensions (i.e. the original two dimensions)

